

Feb. 13, 2002

Date

Express Mail Label No.: EL846174000US

*Don Sampson*

1        THREAD BASED SCALABLE ROUTING FOR AN ACTIVE ROUTER

3                    RELATED APPLICATIONS AND PRIORITY CLAIM

4                    This application is related to prior provisional application no.  
5        60/269,149, and filed February 15, 2001. Applicants claim priority under 35  
6        U.S.C. §119 from that related application.

7                    FIELD OF THE INVENTION

8                    The present invention generally concerns data communications. The  
9        method of the invention more specifically concerns active routers.

10                    BACKGROUND OF THE INVENTION

11                    Active networks offer the promise to improve bandwidth utilization  
12        compared to conventional packet routers, and active routers used in these active  
13        networks can be programmed on a per connection or even a per packet basis. This  
14        programmability makes the routers very flexible, because they are capable of  
15        allocating their finite bandwidth and processing in an application specific manner.  
16        New software applications, which contain protocols developed after a router is  
17        deployed, are able to inject a code for implementing a bandwidth allocation policy  
18        that is best suited for traffic into the network.

19                    Others have investigated providing programmable services with Fast  
20        Ethernet line speeds that implement the NodeOS interface having an active router

1 for programmable execution environments. Generally, the paradigm shift is a goal  
2 to route traffic produced at IEEE 802.3z (gigabit) line speeds with remotely  
3 injected services. However, there are some tradeoffs, such as flexibility and  
4 performance, in the prior methods. For example, a programmable router, such as  
5 Cisco's 7500 line of routers with VIP-4 line cards, offers such routing capacity.  
6 Although the line speeds are similar in Cisco's routers, the VIP-4 processors are  
7 not capable of accepting and then executing remotely injected code while the  
8 router continues to operate.

9 The flexibility versus performance tradeoffs necessarily favor a more  
10 efficient, multiprocessing execution environment. A giga-bit line speed leaves a  
11 window of approximately 8K instructions to process 1 KB packets on a utilized,  
12 dual-issue 500Mhz processor. The ability to efficiently change resource  
13 scheduling is also a great concern, since a 1KB packet arrives every  $7.6\mu s$ , and  
14 may require preempting the packet currently being processed. These tradeoffs  
15 between limiting the amount of processing per packet and the amount of state  
16 (e.g., the amount of data being kept track of per connection basis) the router is  
17 expected to maintain without significantly compromising observed router  
18 flexibility are difficult to balance.

19 Accordingly, there is a need for an improved scalable system routing  
20 method for a gigabit active router, which accounts for the aforementioned  
21 difficulties. There is a further need for an improved routing method, which  
22 reduces system overhead, making it possible to process the packets produced by  
23 gigabit networks.

## 24 SUMMARY OF THE INVENTION

25 These and other needs are met or exceeded by the present scalable  
26 system routing method. Received packets are associating with threads for  
27 processing the received packets, and while processing a previously received  
28 packet, the arrival of an interrupt is checked. If there is an interrupt, a thread is

1 created associating the interrupt is created. Then, a determination of whether the  
2 thread associated with the interrupt has a priority that is higher than the priority of  
3 a thread associated with the previously received packet is made. If the thread  
4 associated with the interrupt has a higher priority than the previously received  
5 packet, the thread associated with the previously received packet is saved in a  
6 Shared Arena storage area. However, if the thread associated with the interrupt  
7 does not have a higher priority than the previously received packet, the thread  
8 associated with the interrupt is queued. Because threads are associated with the  
9 packets, the threads themselves can now be suspended and resumed without  
10 having to disable interrupts, which includes periods during a context switch. As a  
11 result, a more flexible and efficient scheduling routing method can be  
12 implemented.

#### 13 BRIEF DESCRIPTION OF DRAWINGS

14 FIG. 1 illustrates a preferred overall schematic diagram of the  
15 present invention;

16 FIG. 2 illustrates a preferred preemption safe interrupt handler  
17 shown in FIG. 1;

18 FIG. 3 illustrates a preferred interrupt preemptable code for  
19 resuming a kernel thread shown in FIG. 2;

20 FIG. 4 illustrates an overall schematic diagram of the interrupt  
21 handling and the packet processing shown in FIG. 1;

22 FIG. 5 illustrates a preferred nonblocking priority-based run queue;

23 FIG. 6 illustrates an overall diagram of the organization of the kernel  
24 router operating system;

25 FIG. 7 illustrates an unified kernel path implemented with the  
26 present invention is shown in FIG. 7; and,

27 FIG. 8 illustrates a preferred commodity hardware architecture in  
28 which the present invention can be implemented is shown in FIG. 8.

## DETAILED DESCRIPTION OF THE INVENTION

In the present scalable system routing method, received packets are associating with threads for processing the received packets. While a previously received packet is processing, the arrival of an interrupt is checked. If there is an interrupt, a thread is created associating the interrupt is created. Then, a determination of whether the thread associated with the interrupt has a priority that is higher than the priority of a thread associated with the previously received packet is made. If the thread associated with the interrupt has a higher priority than the previously received packet, the thread associated with the previously received packet is saved in a Shared Arena storage area. However, if the thread associated with the interrupt does not have a higher priority than the previously received packet, the thread associated with the interrupt is queued. Because threads are associated with the packets, the threads can now be suspended and resumed without having to disable interrupts, which includes periods during a context switch. As a result, a more flexible and efficient scheduling routing method can be implemented.

## PREFERRED INFRASTRUCTURE OF THE PRESENT INVENTION

The preferred overall schematic diagram of the present invention is shown in FIG. 1, and indicated generally at 10. In the present invention, a packet process 12 is configured to create a thread for each received packet or a group of received packets, which is then managed and routed to the resources according to their priority 14. Because thread creation is sufficiently inexpensive, combining a thread with the processing of each packet is very practical and is preferred. In addition, by creating a thread for each packet, the scheduling of packet processing is now simplified to the scheduling of threads to the processing resources. Thus, the problem of the router providing programmable functionality and routing without specialized hardware is achieved by the present invention. The allocation

of buffer memory is also simplified with the use of threads. Thus, the creation of threads for every packet is an important aspect of the preferred embodiment.

A preferred embodiment assumes a Linux implementation, using kernel threads. In the preferred embodiment, with the creation of kernel threads for every received packets the overall structure of the present active router is preferably organized into two layers, specially an interrupt priority level ("IPL") 16 and a kernel priority level ("KPL") 18, for routing the kernel threads. The packet processing 12 is preferably handled at the kernel priority level 18. The packet process is implemented for creating the kernel threads for the received packets, which is then routed to a plurality of resources. However, when an interrupt is received at the KPL, the packet process also defines the interrupt handlers 20 for handling the interrupt. The interrupt handlers 20 are then executed at the IPL 16. If there is a resume action from the IPL 16, a restartable synchronization 22 of the thread would be handled at the KPL.

However, the present router would not queue or discard a packet interrupt when it is in the process of providing fairness by switching from one thread packet processing to another. Unlike the normal operating systems, which have to protect context switching by disabling interrupts, a Shared Arena 24 is used instead for communication between the two layers 16, 18. In order to use the Shared Arena 24, the interrupt handlers 20 are modified to do interrupt preemptions rather than saves. The difference is in saving the processor's state (e.g., data kept inside the Central Processing Unit) for the interrupted packet processing in the globally accessible Shared Arena, instead of a normal kernel thread private save area in a typical router. In other words, the interrupted thread state is saved 26 in the Shared Arena 24. As a result, the Shared Arena 24 is a communication mechanism that helps to reduce system overhead in taking interrupts, which can eliminate the need to disable interrupts. This preemption makes the interrupted packet processing state restartable by the kernel scheduler on a different processor, and the Shared Arena 24 communicate the

1 synchronization state 28 to the KPL, which executes the restartable  
2 synchronization 22. Thus, as shown in FIG. 1, the invention is preferably  
3 implemented on a multiprocessing hardware 30. However, since the router is now  
4 preempted by an interrupt, rather than disabling the interrupt, the operations (e.g.,  
5 kernel thread resume, save, and critical data structures like the run queue) must be  
6 made preemption safe.

## 7 PREEMPTION SAFE INTERRUPT HANDLER

8 A preferred preemption safe interrupt handler is shown in FIG. 2,  
9 and generally indicated 32. An operation that is preemption safe can be preempted  
10 at any time during its operation and be resumed without corrupting the operation.  
11 However, the procedures that reload the preempted processor state must be made  
12 such that it is restartable, since the procedures too can be preempted. A complete  
13 router based on these operations will improve the handling of high interrupt  
14 (bandwidth) rates. Since the interrupt handlers 20 in the present invention do not  
15 block kernel synchronization, system efficiency does not degrade. The preemption  
16 safe system services will further improve system efficiency by eliminating the need  
17 to disable and reenable interrupts in order to execute these services. Optimizing  
18 these services is essential for good routing performance because the router makes  
19 more use of system scheduling by using suspend and resume to maintain fair  
20 packet processing allocation when there is enough packets (threads) to consider.

21 In addition, since the Shared Arena 24 is accessible in both the KPL  
22 mode 18 and the IPL mode 16, the Shared Arena is not memory that can be  
23 swapped. At the lowest level, there are IPL routines that are invoked directly by  
24 the hardware for external events, such as a timer expiring or a packet arriving from  
25 the network. Typically, an interrupt handler 20 saves the interrupted kernel  
26 thread's 14 context on its kernel stack. The handler 20 then starts executing the  
27 service routine, borrowing some of the interrupted kernel thread's stack. Then,  
28 during the servicing of the interrupt, further interrupts of that type are disabled to

1 prevent stack corruption, overruns, or clobbering the saved state. However, in the  
 2 present invention, the interrupt handler is modified to perform an interruptible  
 3 context switch between the thread 2 that was previously running and an interrupt  
 4 service thread.

5 When there is an interrupt preemption (i.e., interrupt kernel thread),  
 6 thread 2 is first suspended and partially switched. A partial switch is preferably  
 7 performed in order to avoid full thread state being saved to or reloaded from the  
 8 Shared Arena. After thread 2 has been partially switched, the interrupt is handled.  
 9 After the interrupt is complete, thread 2 is only partially resumed, or it is  
 10 completely saved and a higher priority kernel thread is resumed. As a result, this  
 11 save allows the kernel to restart interrupted threads without expecting the  
 12 interrupted kernel thread to first voluntarily yield the resource as it does in the  
 13 prior art, such as the Linux operating system. However, this functionality requires  
 14 a redesign of how and where threads are context switched, saved, and resumed.

#### 15 INTERRUPT PREEMPTABLE CODE

16 A preferred interrupt preemptable code for resuming a kernel thread  
 17 is shown in FIG. 3, and indicated generally at 40. As shown in the previous FIG.  
 18 2, the save areas in the Shared Arena are used as staging areas for the context  
 19 switch. The context switch is used by another thread at KPL to resume a kernel  
 20 thread preempted at IPL. The complexity is in correctly handling the race between  
 21 declaring the thread as having been resumed via setting a shared variable, and an  
 22 another preemption and save that could corrupt the save area. In order to avoid the  
 23 save area from being corrupted, the declaration is setting the identifier of the  
 24 thread that is currently being executed by a resource. This identifier must be set  
 25 before the thread is completely resumed, because the context switch is completed  
 26 by loading and setting the program counter to start executing the code where an  
 27 interrupt asynchronously preempted the kernel thread.

1 In step 1, load directly settable registers (e.g., the working set/data of  
2 the thread that is currently using the processor) from the Share Arena 24, then test  
3 for a nested "restarted" resume in step 2. The next step is to save the indirect  
4 registers in the kernel stack in the Save Arena. After the indirect registers are  
5 saved, the kernel thread can be safely resumed in the process (step 4), and the  
6 process continues by popping remaining registers from the stack in Step 5. The  
7 way context switch, which is safely preempted, is described in cases for the  
8 possible orderings of the race between preemption and declaration.

## 9 FIRST CASE

10 The first case is preemption occurring after declaration but before  
11 the resume is complete. All of the state that has not been loaded into the processor  
12 after the declaration (at the very least the program counter) reside in a memory  
13 location that can be overwritten by another, nested invocation of preemption and  
14 resume. A nested resume occurs after another thread attempts to resume the  
15 context of a thread preempted in the last stages of resume. The nested resume is  
16 detected by doing a check of the preempted program counter in step 2. If the  
17 check returns true, the correct registers have been saved on the stack of the thread  
18 to be resumed, and thus step 3 is skipped. If the check returns false, the correct  
19 registers reside in the save area in the Shared Arena. Since it is possible to  
20 determine the location of the correct register contents after IPL preemption, it is  
21 not necessary to block interrupts for this case.

## 22 THE SECOND CASE

23 The second case is a preemption before declaration that preempts a  
24 thread in the process of yielding its scheduled resource. Conceptually, this  
25 preemption is still equivalent to preempting the yielding thread. However, the  
26 yielding thread at some previous time, acquired mutual exclusion for the resumed  
27 thread's preempted context (through clearing the context's available status in the

1 save area). Suspending the processor at this point would hold up two threads.  
2 Either the yielding context is completely saved in the time to do a voluntary switch  
3 or an interrupt preemption. Therefore, the nonblocking property of the operating  
4 system is maintained without disabling interrupts. The remaining issue is  
5 preventing a voluntary and an interrupt preemption saving to the same memory  
6 location. A voluntary save does not save into the Shared Arena but the thread  
7 descriptor (not shared) as is normally done in multithreading systems. Instead, the  
8 voluntary save completes the necessary code paths in order to provide preemptable  
9 context switching. A faster context switch increases the rate at which the OS is  
10 able to produce simultaneous access to core data structures, principally the run  
11 queue.

## 12 INTERRUPT HANDLING AND PACKET PROCESSING

13 An overall schematic diagram of the interrupt handling and the  
14 packet processing is shown in FIG. 4. In the present invention, the storage for the  
15 packet is allocated on the stack of a kernel thread. Because thread creation is  
16 sufficiently inexpensive, combining a thread with the processing of each packet is  
17 practical. However, by creating a thread for each packet, the scheduling of packet  
18 processing is now simplified to the scheduling of threads to the processing  
19 resources. Thus, the problem of the OS providing direct support, without  
20 specialized hardware operations to change the packet a processor is handling, is  
21 achieved by the present invention. In addition, the allocation of buffer memory is  
22 also simplified.

23 Since the thread descriptor, packet buffer, and thread stack are  
24 allocated as a single object, the interrupt service routine (ISR) is completed by  
25 enqueueing the packet thread to the system run queue. In contrast, direct access to  
26 the run queue from ISR is generally prevented by the run queue's lock in the OS  
27 (e.g., Linux). Once a processing resource becomes available due to a priority  
28 change or a thread exiting, the scheduler starts the thread executing. Since packet

processing is not atomic (i.e., preemptable), low overhead to multiplex the processing of packets on the CPUs for fairness is, thus, allowed. When the thread completes the service initialized to carry out, it enqueues the packet for hard transmit by an output scheduler thread as is the case in Linux. However, unlike in Linux, the output queue can be dynamically varied in length and does not require a lock to protect parallel access. Parallel access on the key data structures that are in the critical path of a packet's trip can help improve the router's utilization of additional routing processors. The combination of preemptable system services and parallel access data structures provide the scalability in the present invention.

## NONBLOCKING PRIORITY SCHEDULING

A preferred nonblocking priority-based run queue implemented with the present invention is shown in FIG. 5. The run queue is accessible for simultaneous parallel accesses by pushing the necessary synchronization into hardware supported atomic operations, including Fetch&Add, Fetch&Swap, Compare&Swap. This technique also makes the run queue accessible to interrupt service routines to make threads runnable without explicitly addressing priority inversion. In open source operating systems (e.g., Linux) and commercial operating systems (e.g., IRIX), a thread lock and a queue lock are implemented to protect the run queue and context switch, respectively. Introducing a parallel access queue would be of questionable use in these operating systems, because locking and its overhead are still required for the context switch. The context switch could race with simultaneous scheduler operations, such as suspend and resume, making a lock to protect thread state necessary. Moreover, a parallel access run queue is less functional than its locked counterpart, since it is unable to service requests for searches and modifications on a snapshot of the run queue. It is assumed that these reasons combined with the complexity of creating a nonblocking operating system have prevented their incorporation in the scheduling core of mainstream operating systems. However, there is a need for handling a

1 faster rate of resource allocation changes than the time slice of these operating  
2 systems (e.g., 10 $\mu$ s to 100 $\mu$ s). There is a possibility that a higher priority packet  
3 will arrive at a line card's streaming rate (7.6 $\mu$ s 1K packets and faster for smaller  
4 packets). As a result, the need for preemptive priorities precludes most existing  
5 nonblocking data structures.

6 However, as shown in FIG. 5, a priority bit vector is provided in a  
7 run queue of the present invention that includes an array of nonblocking Last-In-  
8 First-Out ("LIFO") or First-In-First-Out ("FIFO") data structures. The Top of  
9 Stack ("TOS") is used by a nonblocking LIFO that updates an "Age value" and a  
10 TOS pointer with an atomic operation to add or remove a thread. The "Age value"  
11 prevents a preempted thread from corrupting the queue, since it stores values read  
12 from the queue for a long period of time. However, the queue could be corrupted  
13 by a dequeue, if only the TOS value and next pointer were read before preemption.  
14 Thus, during preemption, the stack is modified, but the TOS value happens to  
15 equal the previously read value when the dequeue is restarted. The dequeue would  
16 atomically swap the previously read TOS value with the wrong next pointer. As a  
17 result, the "Age value" substantially reduces the likelihood of this occurrence.

18 Existing nonblocking LIFO and FIFO algorithms are modified to  
19 indicate at their completion whether a dequeue removed the last entry in a priority  
20 or an enqueue added to a previously empty priority. These are the only two cases  
21 that atomically invert a bit corresponding to the priority emptied or newly  
22 populated, respectively. A single load of this bit vector is used by dequeuers to  
23 locate the queues containing threads. Therefore, enqueue and dequeue operations  
24 normally consume the overhead of a single atomic operation on a loaded system.  
25 The sole caveat is in the unloaded state. If a thread is preempted before  
26 completing the atomic update to the priority bit vector, dequeue attempts can be  
27 directed to an empty priority or a priority with runnable threads, which will be  
28 hidden from dequeue attempts. The first case has to be provided for in the  
29 dequeue algorithm. Simultaneous dequeue attempts can be directed to the same

1 priority. One of the dequeuers masks its local copy of the priority bit vector after it  
2 determined that the priority is really empty. The second case can only be corrected  
3 by periodically sweeping for hidden threads.

#### 4 KERNEL ROUTER OPERATING SYSTEM ORGANIZATION

5 An overall diagram of the organization of the kernel router operating  
6 system is shown in FIG. 6. More specifically, the preferred entry path of the  
7 packets through the operating system is shown. The ideal case that is available in  
8 some commodity hardware gigabit ethernet cards is illustrated. Performing  
9 classification on a packet before buffer allocation presumes that a packet's header  
10 can be efficiently sent before the driver requests that the packet's data be copied.  
11 By placing a classification directly in the interrupt handler, the earliest opportunity  
12 to determine whether a packet should be selectively denied is provided in the  
13 system. The preferred table used in this stage is smaller and faster than the normal  
14 IP lookup table, it is assumed the router services can be applied on logical  
15 groupings of connections rather than individual the source, destination, and TOS  
16 fields.

17 The kernel thread stack, descriptor, and packet-buffer (1500 bytes  
18 MTU) are stored in a single page of memory (4k). Thread allocation uses a private  
19 pool of threads reserved for the router, so reinitialization for each thread is  
20 minimal (i.e. the descriptor of a recently deceased packet processing thread  
21 requires very few updates to safely make it runnable). After the packet is copied to  
22 the buffer, the thread is added to the global run queue. The processors, then, poll  
23 the run queue for new threads after a thread exits, a contended synchronization  
24 operation, a timer interrupt event, or constantly by the lowest priority idle thread.  
25 The packet processing thread is initialized to start executing at the address for the  
26 packet processing routine. The packet processes to termination unless it  
27 voluntarily yields the resource. Once packet processing is completed, it is  
28 scheduled for an output interface through a nonblocking output queue.

1           The nonblocking output buffer is similar in design, which includes a  
2 hybrid private and parallel access run queue. Since hardware supported atomic  
3 operations are more time costly, they should be avoided unless contention  
4 mandates their use. The output queue is optimized for parallel enqueue from  
5 threads submitting completed packets. The timer interrupt only occurs on one  
6 CPU because a single thread is used to service the timer interrupt. This  
7 optimization enables the use of normal operations, rather than atomic operations to  
8 dequeue the packets. However, parallel dequeue is not as useful because the  
9 system bus serializes their transmission to the transmission buffer.

## 10 AN EXECUTION MODEL FOR A SCALABLE ROUTER

11           An unified kernel path implemented with the present invention is  
12 shown in FIG. 7. In the present invention, the programmability is preferably  
13 moved down to the level of a trusted path within the kernel. In particular, this is  
14 the level at which the Core Network Context (CNC) and User Network Contexts  
15 (UNCS) hierarchy are implemented from the Protean architecture invented by R.  
16 Sivakurnar, S. Han, and V. Bharghavan disclosed in an article entitled, "A scalable  
17 architecture for active networks," published in OpenArch, 2000. The trusted path  
18 first executes CNC services that provide the policies applied to all traffic through  
19 the router. These policies, in turn, provide fair allocation of output bandwidth and  
20 processing capacity. The CNC manages a predefined number of UNCS, so the  
21 overhead of providing fairness and scheduling is reduced in operations like  
22 classifying amongst a limited number of UNCS. UNCS are an abstraction for  
23 differentiating services based on routing source, destination, and type of service.  
24 The logical grouping of connection types into UNCS places a tractable limit on the  
25 amount of state maintained in the router. This architecture has demonstrated the  
26 benefit of providing programmable services comparable to more flexible active  
27 network designs.

1           The trusted path can provide a way to remotely inject into a router  
2 both programmable services and execution environments for untrusted services.  
3 Some execution environments, such as the Java interpreter for ANTS, require  
4 direct access to system services such as thread management. Since services are  
5 executed directly in the kernel, an execution environment for untrusted services  
6 can be accomplished without significant overhead.

## 7 COMMODITY HARDWARE ARCHITECTURE

8           A preferred commodity hardware architecture in which the present  
9 invention can be implemented is shown in FIG. 8. The preferred hardware is a two  
10 stage architecture. The lower tier is composed of commodity PIII 550Mhz systems  
11 (not shown) that act as the router's workhorse. This is where the routing  
12 operations for an input packet are executed. The theoretical capacity of the  
13 processors at these nodes is 17.6 cycles per word produced by the external gigabit  
14 ethernet interface. It is necessary to attain more processing capacity, since the real  
15 capacity is far lower due to memory and bus overheads. It is less expensive to add  
16 processing capacity with a second processor than increasing the processor speed.  
17 In practice, it is expected that the capacity will be roughly doubled without static  
18 partitioning. The design of the present invention helps to remove the Operating  
19 System as a scaling bottleneck and provides scalable system services for  
20 downloaded drivers.

21           Once processed control information is appended to the packet for  
22 routing within the router architecture, there may still be queueing issues between  
23 the programmable nodes of the router. Thus, the output packet scheduler  
24 implements guarantees for Quality of Service ("QoS") and rate control on  
25 forwarding packets within the router. Although the design shown is limited to  
26 controlling bandwidth, the present invention can be extended to handle other  
27 controls, such as latency and jitter. As a result, these various other  
28 implementations are contemplated and are within the scope of the present

1 invention. Another enhancement considered is internal support for multicast to the  
2 routing nodes. One of the goals of the present invention is to provide gigabit  
3 routing services at the lower tier with inexpensive multiprocessors and one  
4 external and one internal gigabit interface. In the upper tier, a commodity gigabit  
5 switch acts as a low latency router fabric, connecting the lower tier.

6 While various embodiments of the present invention have been  
7 shown and described, it should be understood that other modifications,  
8 substitutions and alternatives are apparent to one of ordinary skill in the art. Such  
9 modifications, substitutions and alternatives can be made without departing from  
10 the spirit and scope of the invention, which should be determined from the  
11 appended claims.

12 Various features of the invention are set forth in the appended  
13 claims.